

Є.В. Івохін¹, Г.В. Шелякін², М.Ф. Махно³^{1,2,3}Київський національний університет імені Тараса Шевченка, Україна
пр. Глушкова, 4д, м. Київ, 03680¹ivohin@knu.ua²shelyakingleb17@gmail.com³makhnom@gmail.com¹<https://orcid.org/0000-0002-5826-7408>²<https://orcid.org/0009-0002-7171-6535>³<https://orcid.org/0000-0001-9694-6602>

УДОСКОНАЛЕННЯ МЕТОДУ КОЛАБОРАТИВНОЇ ФІЛЬТРАЦІЇ ШЛЯХОМ ІНТЕГРУВАННЯ СЕМАНТИЧНОГО ТА ЧАСОВОГО ФАКТОРІВ І МЕТОДУ КЛАСТЕРНОГО АНАЛІЗУ

Анотація. У статті розглядається алгоритм формування рекомендацій на основі колаборативної фільтрації з урахуванням впливу семантичного і часового факторів та його удосконалення за допомогою методів кластерного аналізу з метою зменшення навантаження на рекомендаційну систему та покращення якості рекомендацій шляхом відсіювання незмістовного контенту і збереження контексту під час генерації рекомендацій. Проаналізовано вплив семантичного та часового факторів на якість роботи системи рекомендацій (похибка при апроксимації оцінки) та застосування методу кластерного аналізу на швидкодії системи при великому наборі даних. Запропоновано методику прискорення обробки отриманих даних про користувачів, яка полягає у спробі врахувати факт зміни інтересів користувачів з часом і можливість розбити контент статистичних даних за сукупністю конкретних ознак. Сформульовано процедуру попередньої обробки даних (агрегація даних) для методу колаборативної фільтрації на основі порівнянь об'єктів з використанням методу кластеризації, що дозволило зменшити складність обчислень і, відповідно, час формування рекомендацій. Наведено алгоритм підрахунку оцінки об'єкта з урахуванням часового та семантичного факторів. Розроблено програмне забезпечення, перевірено адекватність роботи запропонованого методу, використовуючи набори даних з різних доменних областей. У результаті перевірки було виявлено, що модифікований алгоритм має кращі показники ефективності роботи у порівнянні з найвним методом.

Ключові слова: рекомендаційні системи, колаборативна фільтрація, кластерний аналіз, метод семантичної подібності, часовий фактор.

Ye. Ivohin¹, G. Shelyakin², M. Makhno³^{1,2,3}Taras Shevchenko National University of Kyiv, Ukraine
Glushkova st., 4d, Kyiv, 03680¹ivohin@knu.ua²shelyakingleb17@gmail.com³makhnom@gmail.com¹<https://orcid.org/0000-0002-5826-7408>²<https://orcid.org/0009-0002-7171-6535>³<https://orcid.org/0000-0001-9694-6602>

IMPROVING THE METHORD OF COLLABORATIVE FILTERING BY INTEGRATING SEMANTIC AND TEMPORAL FACTORS AND THE METHORD OF CLUSTER ANALYSIS

Abstract. The article examines the algorithm for generating recommendations based on collaborative filtering, taking into account the influence of semantic and time factors and its improvement using cluster analysis methods in order to reduce the load on the recommendation system and improve the quality of recommendations by filtering out meaningless content and preserving the context during the generation of recommendations. The impact of semantic and time factors on the quality of the recommendation system (error in estimation approximation) and the application of the cluster analysis method on the speed of the system with a large set of data are analyzed. A technique for accelerating the processing of received data about users is proposed, which consists in an attempt to take into account the fact that users' interests change over time and the possibility of breaking down the content of statistical data by a set of specific features. A data preprocessing procedure (data aggregation) was formulated for the method of collaborative filtering based on comparisons of objects using the clustering method, which made it possible to reduce the complexity of calculations and,

accordingly, the time for the formation of recommendations. An algorithm for calculating the object's assessment is presented, taking into account temporal and semantic factors. The software was developed, the adequacy of the proposed method was verified using data sets from different domain areas. As a result of the verification, it was found that the modified algorithm has better performance indicators compared to the naive method.

Keywords: recommender systems, collaborative filtering method, cluster analysis, semantic similarity method, time factor.

Вступ

У інформаційному суспільстві, де кількість доступної інформації зростає з неймовірною швидкістю, проблема "перевантаження інформацією" стає все більш актуальною. Люди щодня стикаються з неперервним потоком новин, соціальних медіа-повідомлень, електронних листів, реклами, а також з величезною кількістю контенту на стрімінгових платформах і в електронних магазинах. Це створює інформаційний шум, який може паралізувати здатність людини до прийняття рішень. В такому середовищі системи рекомендацій відіграють ключову роль, допомагаючи фільтрувати непотрібну інформацію та визначати те, що найбільш релевантно для конкретного користувача.

Існує багато підходів для створення рекомендаційних систем, серед яких потрібно відзначити, наприклад, ті, що базуються на використанні нейронних мереж або статистичних математичних моделей. Прикладом таких підходів є методика залучення алгоритму колаборативної фільтрації на основі інформації про користувача (user-based) та на основі порівнянь об'єктів рекомендацій (item-based). Проте, з урахуванням зростання кількості контенту та користувачів обчислювальні потужності серверів можуть бути недостатніми для ефективної роботи вищезазначеного методу. В даному дослідженні пропонується враховувати той факт, що інтереси користувачів можуть змінюватись з часом, та те, що контент можна заздалегідь розбити на підмножини за певними ознаками (класифікувати або кластеризувати), що дозволить рекомендаційній системі швидше обробляти дані. В рамках проведеного дослідження зроблено спробу удосконалити метод колаборативної фільтрації на основі порівняння об'єктів із використанням часового фактора і семантичної

подібності та завдяки застосуванню методики кластерного аналізу.

Актуальність

Актуальність тематики дослідження підтверджується підвищеною увагою з боку сучасних науковців (див., наприклад, [1-4]). Серед найбільш відомих публікацій необхідно виділити статтю С. Бірковського, І. Сантадора та Д. Тікка [1], яка присвячена ґрунтовному дослідженню методики та алгоритмів формування рекомендацій контенту. В роботі розглянуто можливі виклики, з якими можна зіткнутися під час розробки рекомендаційних систем, і запропоновано практичне застосування методів оптимізації. В роботі А. Хінебурга та А.Д. Кейма [2] наведено новий алгоритм на основі оцінки щільності ядра для кластеризації у великих мультимедійних базах даних, в якому кластери можна ідентифікувати шляхом визначення атракторів щільності, а кластери довільної форми можна легко описати простим рівнянням загальної функції щільності. Ж. Сандер у своїй роботі [3] узагальнив відомий алгоритм DBSCAN (отримав назву GDBSCAN), який дозволяє кластеризувати точкові об'єкти, а також просторово розширені об'єкти відповідно до їхніх просторових і непросторових атрибутів. Крім того, ним було представлено чотири програми, що використовують 2D-точки (астрономія), 3D-точки (біологія), 5D-точки (науки про Землю) і 2D-полігони (географія), які демонструють можливість застосування GDBSCAN до вирішення реальних проблем.

Постановка задачі

Розглянемо процедуру прийняття (вибору) деякого рішення з урахуванням наявної сукупності рекомендацій. Припустимо, що об'єктами рекомендацій будуть деякі сутності, характеристики яких можна розділити на три групи:

1. Оцінки. Під оцінкою розуміється вектор відповідності між користувачами та об'єктом, кожна з яких визначається у вигляді деякого числового значення. Передбачається, що відповідне значення є обмеженим знизу та згори, наприклад, від 1 до 5, від 0 до 1 тощо.

2. Текстова інформація. Текстовою інформацією вважається деякий скінченний список характеристик кожного об'єкта. До елементів текстової інформації відносяться, наприклад, описи, назви, резюме тощо.

3. Часова інформація. Під часовою інформацією розуміють дані про час та тривалість виконання певної дії користувача над деяким об'єктом. Наприклад, часову інформацію можна отримати при перегляді та оцінюванні об'єкта, рецензуванні об'єкта, наданні відгуку щодо об'єкта, замовленні об'єкта тощо.

Тоді задача полягає у тому, щоб розробити алгоритм, за допомогою якого можна без обмежень на проблемну (доменну) область отримувати апроксимації оцінок користувачів з використанням згаданих вище характеристик і який дозволяв робити це швидше, ніж найвний метод колаборативної фільтрації.

Результати

Колаборативна фільтрація – це метод рекомендації контенту, що базується на реакціях користувачів на контент. Вперше метод колаборативної фільтрації був запропонований Голдбергом у 1992 році у вигляді системи фільтрації електронної пошти [4]. З того часу він став предметом розвідок науковців [1, 5-8] та ін.

Головною метою методу є розрахунок оцінки контенту, з яким користувач ще не знайомий, використовуючи інформацію про його попередні реакції. Чим ближче оцінка наближається до реальної за умови зменшення похибки, тим більш якісним буде кінцевий рекомендаційний висновок. Для отримання найбільш об'єктивної

оцінки рекомендаційної системи необхідно мати якомога більше оцінок користувачів і способи їх аналізу. Таким чином, застосування методики колаборативної фільтрації дозволяє зробити більш точні рекомендації користувачам на основі їх власних реакцій на контент.

З наведеного випливає, що основними поняттями, які застосовуються у методиці є користувач, контент, оцінка та рекомендація.

Користувачем є особа, яка має бути зареєстрованою на деякому сервісі та має можливість перегляду, оцінки, покупки контенту на ньому тощо.

Контент (об'єкт) представляє собою вміст сайту, з яким користувач взаємодіє в рамках сервісу та який може бути оцінений користувачами на основі сукупності метаданих: опису, назви, вартості, тегів тощо.

Оцінка визначається числовим значенням з множини натуральних чисел, яке є обмеженим нулем знизу та деяким додатнім цілим числом згори.

Під рекомендацією розуміємо множину, що може складатися з групи змістовних пропозицій (описів), яка буде створена рекомендаційною системою для користувача.

Оцінки користувачів для об'єктів на зареєстрованому сервісі зручно представляти у вигляді матриці $R = \{R_{ui}\}$, $u = \overline{1, n}$, $i = \overline{1, m}$, де n – кількість користувачів сервісу та m – кількість об'єктів.

Розглянемо довільного користувача сервісу, який визначається деяким порядковим номером u , $u = \overline{1, n}$. Зміст нашої задачі полягає у тому, щоб апроксимувати можливу оцінку користувача u щодо об'єкту i , $i = \overline{1, m}$.

Процедуру апроксимації можна описати так:

1. Для кожного об'єкта i розраховуємо, наскільки він схожий з об'єктом i .
2. Формуємо множину об'єктів, які є найбільш схожими до об'єкта i .
3. Розраховуємо оцінку об'єкта i на основі оцінок з множини найбільш схожих об'єктів.

Тепер опишемо більш детально

кожний з наведених кроків даної процедури. Оскільки кожному об'єкту відповідає стовпчик матриці, то будемо розраховувати міру схожості за стовпцями, використовуючи, наприклад, скориговану схожість косинусів. Формула для оцінки схожості $sim(i, \hat{i})$ об'єктів i та \hat{i} може бути записана так:

$$sim(i, \hat{i}) = \frac{\sum_{u=1}^n (R_{iu} - \bar{R}_u)(R_{i\hat{u}} - \bar{R}_u)}{\sqrt{\sum_{u=1}^n (R_{iu} - \bar{R}_u)^2} \sqrt{\sum_{u=1}^n (R_{i\hat{u}} - \bar{R}_u)^2}} \quad (1)$$

де \bar{R}_u – середнє значення оцінок користувача u .

Отримана величина має такі властивості:

$$- sim(i, \hat{i}) \in [0, 1]; \quad - sim(i, \hat{i}) = 0,$$

якщо користувач не поставив оцінки, або дорівнює 1, якщо об'єкти співпали.

Далі обираємо задану кількість об'єктів $s > 0$, які найбільш схожі з обраним об'єктом i . Для цього потрібно відсортувати отримані величини за спаданням та обрати s перших об'єктів. Сортування гарантуватиме те, що перші s об'єктів матимуть найбільші величини схожості, які є близькими до 1.

Перенумеровуючи ці об'єкти, обчислимо, яку можливу оцінку поставив би користувач u об'єкту i . Для цього використаємо формулу:

$$p_{ui} = \frac{\sum_{i=1}^s R_{ui} sim(i, \hat{i})}{\sum_{i=1}^s sim(i, \hat{i})}. \quad (2)$$

Зміст даної оцінки полягає у тому, що чим більша міра схожості об'єктів i та \hat{i} , тим більший вклад об'єкт \hat{i} вносить до остаточної оцінки.

У дослідженні пропонується методика прискорення обробки отриманих про користувачів даних, яка полягає у спробі врахувати факт зміни інтересів користувачів з часом, і можливість розбити контент статистичних даних за певними ознаками.

Запропоновано підхід на основі методу кластерного аналізу та методу

колаборативної фільтрації із забезпеченням порівняння об'єктів з урахуванням часового фактора та семантичної подібності.

Для того, щоб врахувати часовий фактор, використовується функція “старіння” інформації у вигляді $f(t) = e^{-t}$, $0 \leq t \leq \infty$, $f(t) \in [0, 1]$. Ця функція монотонно спадає, зі збільшенням значення t на виході отримуємо менші значення. Алгоритм використання функції старіння має вигляд:

1. Розраховуємо для кожного об'єкта значення функції “старіння” для показника часової інформації.

2. Відсортовуємо отримані значення за спаданням. Це гарантуватиме, що об'єкти, які були переглянуті нещодавно, матимуть більший вплив на рекомендацію.

3. Обираємо перші $v > 0$ об'єктів, які наразі знаходяться у тренді користувача та, відповідно, мають більший вплив на те, що він хоче бачити.

Для врахування семантичного фактора введемо такі поняття: M - множина текстової інформації, що відображає суть деякого об'єкта i . Кожен елемент множини є вектор (a_1, a_2, \dots, a_K) , де a_j - це деяке текстове представлення j -о атрибута поточного об'єкта $j = \overline{1, K}$. Це, як вже було сказано раніше, можуть бути описи, рецензії, набір технічних характеристик, відгуки та інше.

Тоді міру близькості двох об'єктів i та \hat{i} можна представити такою формулою:

$$sim_{meta}(i, \hat{i}) = \frac{\sum_{j=1}^K sim(i_j, \hat{i}_j)}{K} \quad (3)$$

де $sim(i_j, \hat{i}_j)$ - коефіцієнт подібності j -ї характеристики об'єктів i та \hat{i} із застосуванням, наприклад, сіамських нейронних мереж у якості оцінювачів.

Визначимо квадратну симетричну матрицю $S = \{S_{ij}\}$, $i, j = \overline{1, m}$, $S_{ij} = S_{ji}$, $i \neq j$, де m - кількість об'єктів на сервісі, S_{ij} - оцінки схожості i -го та j -го об'єктів. Це дозволяє зберігати результати оцінювання у базі даних, розраховувати оцінки в

офлайн режимі з урахуванням того, що кількість об'єктів зростає доволі повільно. Крім цього, такий підхід позитивно відбивається на швидкості оновлення даних.

З іншого боку, можна сформулювати процедуру попередньої обробки даних (агрегації даних) для застосування методу колаборативної фільтрації на основі порівнянь об'єктів з використанням методу кластеризації, яка полягає в таких кроках:

1. Будується матриця подібності об'єктів з використанням наявних даних у вибірці.

2. Проводиться розбиття множини об'єктів на відповідні кластери на основі алгоритмів кластеризації.

3. Дані про відповідні кластери зберігаються для майбутнього використання.

Агрегація даних перед початком роботи алгоритму може зменшити час формування рекомендацій та складність необхідних для цього обчислень. За умов використання необроблених даних алгоритм формування рекомендацій має виконувати додаткову роботу з обробки та групуванням даних під час виконання. Це, як правило, займає багато часу та потребує відповідних витрат обчислювальних ресурсів.

На основі наведених способів оцінювання подібності можна вивести остаточний алгоритм підрахунку оцінки об'єкта.

Покладемо

$$U(i, \hat{i}) = \alpha \text{sim}_{\text{mark}}(i, \hat{i}), \quad (4)$$

де $\text{sim}_{\text{mark}}(i, \hat{i}) \in [0,1]$ - міра схожості об'єктів за оцінками користувачів, а $\alpha \in [0,1]$ - ваговий коефіцієнт, який відповідає за величину впливу оцінки $\text{sim}_{\text{mark}}(i, \hat{i})$ в остаточній оцінці.

Аналогічно визначимо

$$V(i, \hat{i}) = \beta \text{sim}_{\text{meta}}(i, \hat{i}), \quad (5)$$

де $\text{sim}_{\text{meta}}(i, \hat{i}) \in [0,1]$ - міра схожості об'єктів на основі порівняння їх метаданих, а $\beta \in [0,1]$ - ваговий коефіцієнт, який відповідає за величину впливу оцінки $\text{sim}_{\text{meta}}(i, \hat{i})$ в остаточній оцінці.

Тоді можна записати остаточну

оцінку схожості (близькості) об'єктів i та \hat{i} , використовуючи наведені вище величини:

$$\text{sim}_{\text{hybrid}}(i, \hat{i}) = (U(i, \hat{i}) + V(i, \hat{i})) / 2. \quad (6)$$

Очевидно, що запропонована оцінка подібності повністю співпадає з властивостями оцінки методу колаборативної фільтрації на основі порівнянь об'єктів.

Тепер побудуємо остаточну формулу розрахунку оцінки користувача u для об'єкта i з урахуванням старіння та семантичного фактора:

$$p_{ui} = \frac{\sum_{i=1}^v R_{ui} \text{sim}_{\text{hybrid}}(i, \hat{i}) f(t_i)}{\sum_{i=1}^v \text{sim}_{\text{hybrid}}(i, \hat{i}) f(t_i)}. \quad (7)$$

Для розробки програмного забезпечення було використано мову програмування Python та безкоштовне середовище для написання коду VS Code.

Проведено перевірку адекватності роботи запропонованого методу, використовуючи набори даних з різних проблемних областей, отримано показники швидкості обчислень рекомендацій.

На діаграмі (рис.1) наведено показники середньої швидкості обчислень відповідних апроксимацій (колонка синього кольору вказує на час, який був необхідний для виконання модифікованого алгоритму (300 мс), а помаранчева колонка - на час, який був необхідний для виконання наївного модифікованого алгоритму (500 мс)).

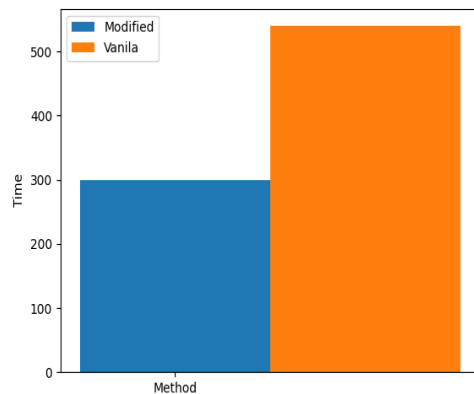


Рис. 1. Порівняння швидкості роботи наївного та модифікованого методів

Ці дані підтверджують, що модифікований алгоритм демонструє більш швидкий режим роботи у порівнянні з наївним методом.

У якості прикладу практичного застосування запропонованого методу розглянемо сервіс з оцінювання фільмів та проаналізуємо можливі при цьому похибки в результатах формування рекомендацій. Отримаємо результати статистичного спостереження процедури оцінювання користувачами випадково обраного об'єкта рекомендацій. Позначимо через y – величину похибки (середньоквадратичне відхилення), x – номер обраного об'єкта. Отримані результати наведено на рис.2 та 3. Колонки синього кольору показують величини похибки між оцінками на основі модифікованого методу, а помаранчеві колонки показують величину похибки для

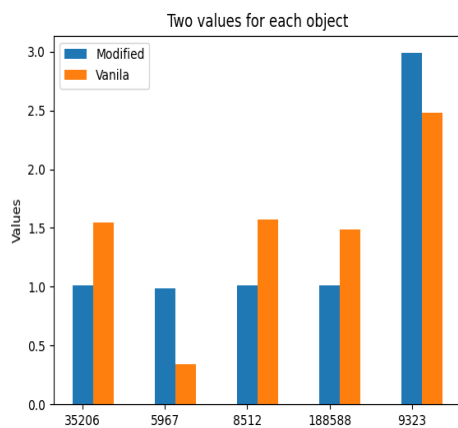


Рис. 2. Порівняння величин похибок між наївним та модифікованим методами

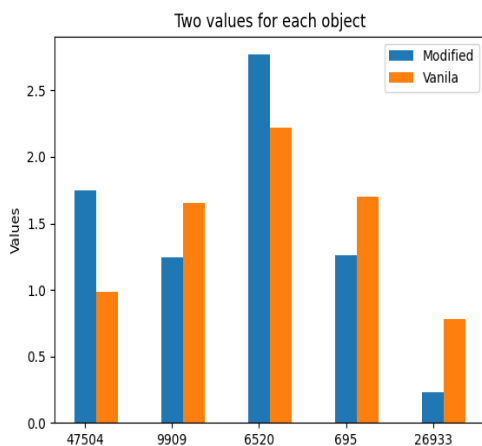


Рис. 3. Порівняння величин похибок між наївним та модифікованим методами

наївного методу. На діаграмах видно, що у більшості випадків похибка модифікованого методу є меншою, ніж похибка наївного методу.

В іншому прикладі (формування рекомендацій у сфері музичних стрічок) за аналогічними отриманими даними результатів статистичного спостереження розраховано результати порівняння похибок для розробленого та наївного методів (рис.4,5). Колонка синього кольору показує величину похибки між оцінками для модифікованого методу, а помаранчева колонка - величину похибки оцінками для наївного методу. Як і раніше, маємо, що у більшості випадків похибка модифікованого методу є меншою порівняно з похибкою наївного методу.

Отримані в проведеному дослідженні результати дають підстави стверджувати

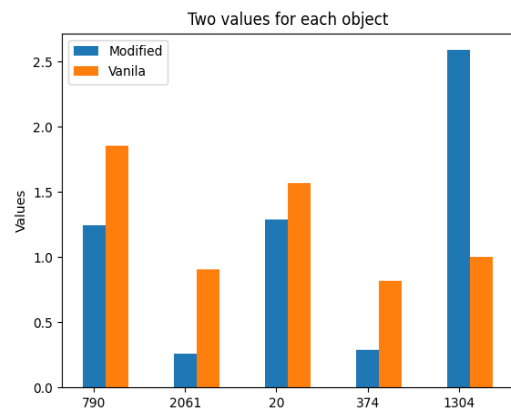


Рис. 4. Порівняння величин похибок між наївним та модифікованим методами

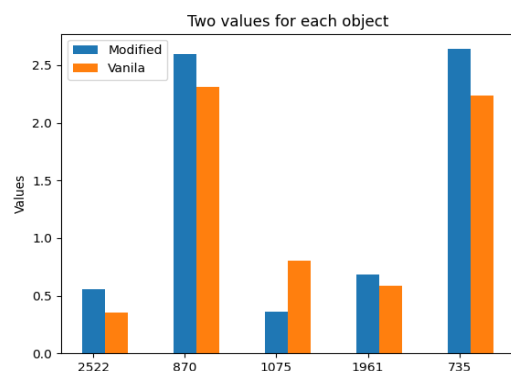


Рис. 5. Порівняння величин похибок між наївним та модифікованим методами

про доцільність використання запропонованої методики, завдяки якій можна покращити швидкість роботи та якість висновків рекомендаційних систем у різних проблемних областях.

Висновки

Запропоновано підхід для удосконалення методу колаборативної фільтрації з використанням часового та семантичних факторів шляхом впровадження спеціальних показників, які характеризують вхідні дані для їх попередньої обробки. Це дало змогу використовувати запропоновану модифікацію у різних проблемних (доменних) областях на основі розробленого програмного забезпечення. Перевірено адекватність роботи запропонованого методу, використовуючи набори даних з різних доменних областей.

У результаті проведених обчислень було встановлено більш високу ефективність роботи розробленого алгоритму формування рекомендацій у порівнянні з наївним методом.

Результати дослідження можуть знайти своє застосування у сфері реалізації різних інтернет-сервісів, серед яких можна виділити сервіси інтернет-речей, стрімінговий сервіс, обслуговування інтернет-магазинів тощо.

Література

1. Berkovsky S., Cantador I., Tikk D. Collaborative recommendations: algorithms, practical challenges and applications. (2019). World scientific publishing. <https://doi.org/10.1142/11131>
2. Hinneburg A., Keim D.A. A general approach to clustering in large databases with noise (2003). Knowledge and Information Systems, 2003. 5 (4). P. 384–415. <https://doi.org/10.1007/s10115-003-0086-9>
3. Sander J. Density-based clustering in Spatial Databases (1998). Data Mining and Knowledge Discovery. Vol.2. P. 169–194.
4. Goldberg D., Nichols D., Oki B.M., Terry D. Using collaborative filtering to weave an information tapestry (1992). Commun ACM. P.61–70. <https://doi.org/10.1145/138859.138867>
5. Horasan F., Yurtttakal A., Gunduz S. A novel model based collaborative filtering recommender system via truncated ULV decomposition (2023). Journal of King Saud University – Computer and Information Sciences. V.35 (8). <https://doi.org/10.1016/j.jksuci.2023.101724>

6. Marappan R. Recommender system for movielens datasets using an item-based collaborative filtering in Python (2022). International Journal of Mathematical, Engineering, Biological and Applied Computing, 1(1). – P. 42–43. <https://doi.org/10.31586/ijmebac.2022.340>
7. Natarajan S., Vairavasundaram S., Natarajan S., Gandomi A.H. Resolving data sparsity and cold start problem in collaborative filtering recommender system using linked open data (2020). Expert Syst Appl. <https://doi.org/10.1016/j.eswa.2020.113248>
8. Shen J., Wei Y., Yang Y. Collaborative filtering recommendation algorithm based on two stages of similarity learning and its optimization. Computers Materials & Continua. 58(2). P. 659–674. <https://doi.org/10.32604/cmc.2019.05858>

References

1. Berkovsky S., Cantador I., Tikk D. Collaborative recommendations: algorithms, practical challenges and applications. (2019). World scientific publishing. <https://doi.org/10.1142/11131>
2. Hinneburg A., Keim D.A. A general approach to clustering in large databases with noise (2003). Knowledge and Information Systems, 2003. 5 (4). P. 384–415. <https://doi.org/10.1007/s10115-003-0086-9>
3. Sander J. Density-based clustering in Spatial Databases (1998). Data Mining and Knowledge Discovery. Vol.2. P. 169–194.
4. Goldberg D., Nichols D., Oki B.M., Terry D. Using collaborative filtering to weave an information tapestry (1992). Commun ACM. P.61–70. <https://doi.org/10.1145/138859.138867>
5. Horasan F., Yurtttakal A., Gunduz S. A novel model based collaborative filtering recommender system via truncated ULV decomposition (2023). Journal of King Saud University – Computer and Information Sciences. V.35 (8). <https://doi.org/10.1016/j.jksuci.2023.101724>
6. Marappan R. Recommender system for movielens datasets using an item-based collaborative filtering in Python (2022). International Journal of Mathematical, Engineering, Biological and Applied Computing, 1(1). – P. 42–43. <https://doi.org/10.31586/ijmebac.2022.340>
7. Natarajan S., Vairavasundaram S., Natarajan S., Gandomi A.H. Resolving data sparsity and cold start problem in collaborative filtering recommender system using linked open data (2020). Expert Syst Appl. <https://doi.org/10.1016/j.eswa.2020.113248>
8. Shen J., Wei Y., Yang Y. Collaborative filtering recommendation algorithm based on two stages of similarity learning and its optimization. Computers Materials & Continua. 58(2). P. 659–674. <https://doi.org/10.32604/cmc.2019.05858>

The article has been sent to the editors 16.01.24.

After processing 13.02.24.

Submitted for printing 20.03.24.

Copyright under license CCBY-SA4.